



Monitoring I/O on Data-Intensive Clusters

Visualizing Disk Reads and Writes on Hadoop MapReduce Jobs

Thursday, July 31

Joel Ornstein



Joshua Long



Carson Wiens



Mentors: Steve Senator, Tim Randles, Vaughan Clinton, Mike Mason, Graham Van Heule – HPC 3



Operated by Los Alamos National Security, LLC for NNSA

UNCLASSIFIED

LA-UR-14-26019

1



Background

Motivation:

- I/O Intensive Jobs
 - Large amounts of scientific data

Background

Motivation:

- I/O Intensive Jobs
 - Large amounts of scientific data

Traditional HPC

- Limiting factor mostly lies in processing speed

Background

Motivation:

- I/O Intensive Jobs
 - Large amounts of scientific data

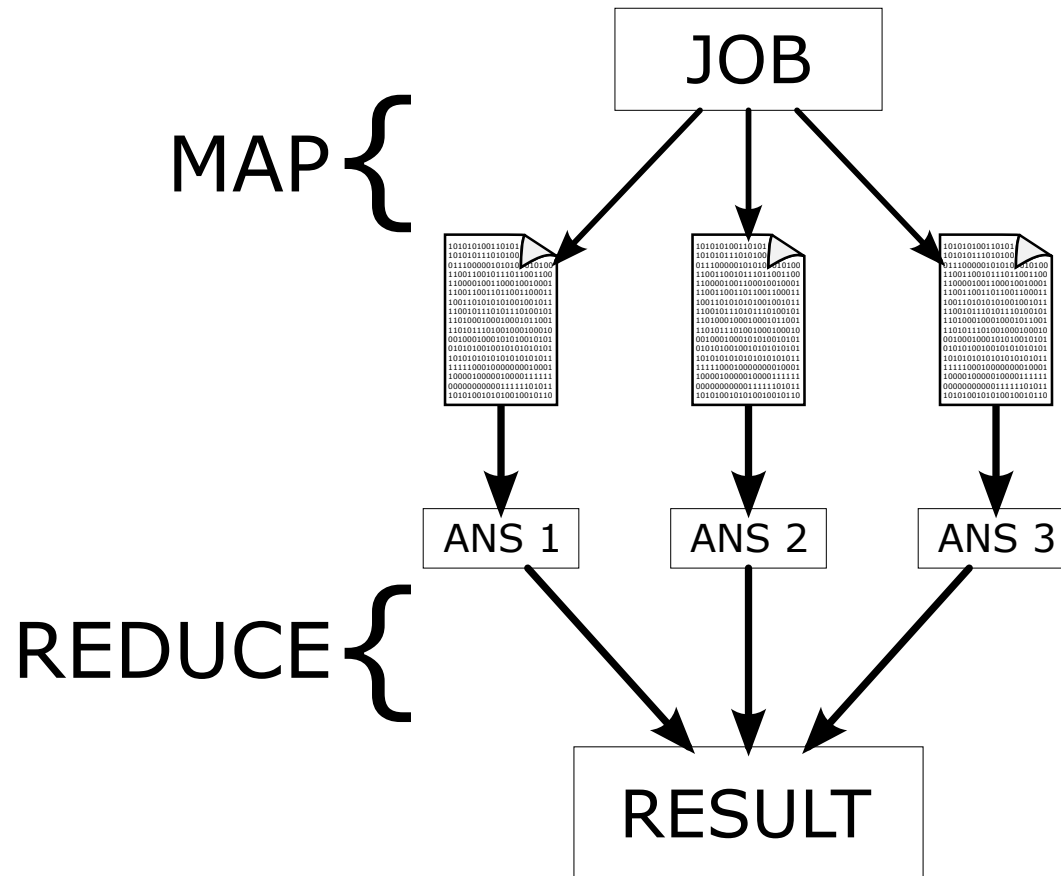
Traditional HPC

- Limiting factor mostly lies in processing speed

I/O Intensive Jobs

- Bottlenecked by read/write disk speed
- MapReduce
 - Move jobs to the data (instead of vice-versa)

MapReduce



I/O Monitoring

Why?

- Nodes break
- Jobs run without using the specified resources

I/O Monitoring

Why?

- Nodes break
- Jobs run without using the specified resources

Deliverables

- Programs that are helpful for monitoring a Hadoop 2.3 cluster
 - Splunk App for HadoopOps
 - Ganglia
 - Other methods

I/O Monitoring

Why?

- Nodes break
- Jobs run without using the specified resources

Deliverables

- Programs that are helpful for monitoring a Hadoop 2.3 cluster
 - Splunk App for HadoopOps
 - Ganglia
 - Other methods
- Data tests
 - bonnie++
 - teragen and terasort

Environment

- 11-node CentOS cluster
 - 1 head node and 10 compute nodes
- FDR InfiniBand 56-Gb/second
 - IP over IB
 - Faster than disks can read/write
- Hadoop 2.3.0
- MRv2/YARN
 - Yet Another Resource Negotiator
 - Runs MapReduce jobs in Hadoop environment
- Java 1.6

Monitoring Tools

Splunk

- software for searching and analyzing logs
- able to generate graphs, charts, gauges, etc.
- web interface

Monitoring Tools

Splunk

- software for searching and analyzing logs
- able to generate graphs, charts, gauges, etc.
- web interface

Ganglia

- software for monitoring clusters
- generates plots from input
- web interface

Monitoring Tools

Splunk

- software for searching and analyzing logs
- able to generate graphs, charts, gauges, etc.
- web interface

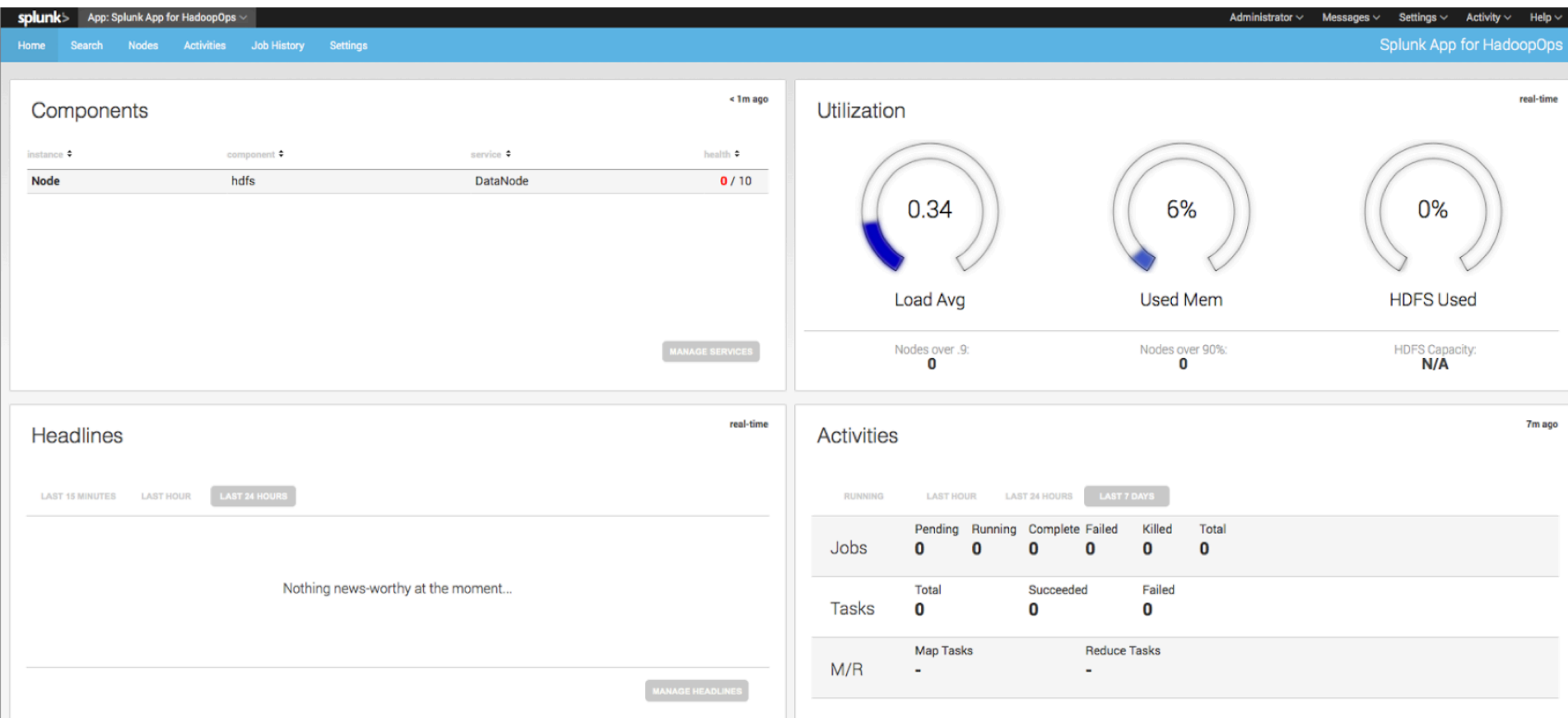
Ganglia

- software for monitoring clusters
- generates plots from input
- web interface

iostat

- outputs I/O statistics for devices
- command-line interface

Splunk App for HadoopOps



Ganglia

**Goldenrod Cluster Cluster Report for Tue, 22 Jul 2014 09:34:31 -0600**

Get Fresh Data

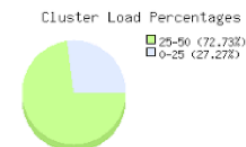
Metric Last Sorted

Physical View

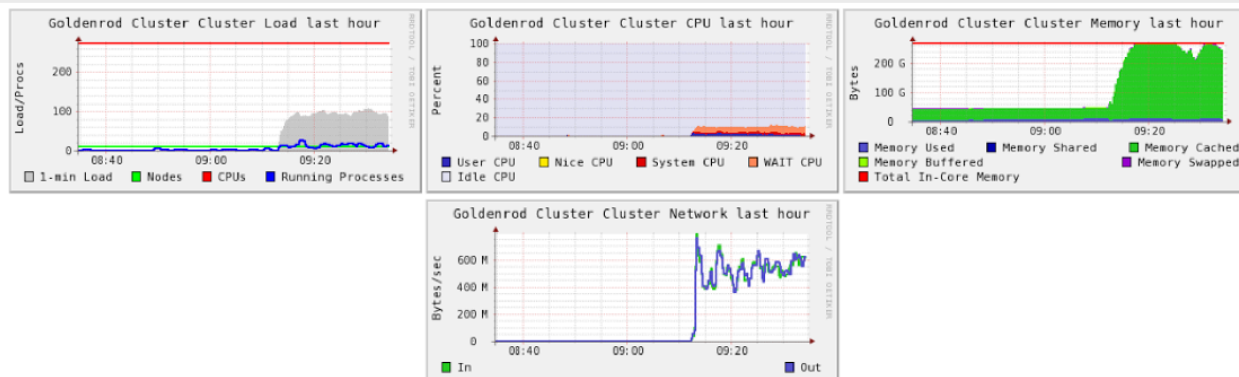
Grid > Goldenrod Cluster >

CPU's Total: 272
Hosts up: 11
Hosts down: 0

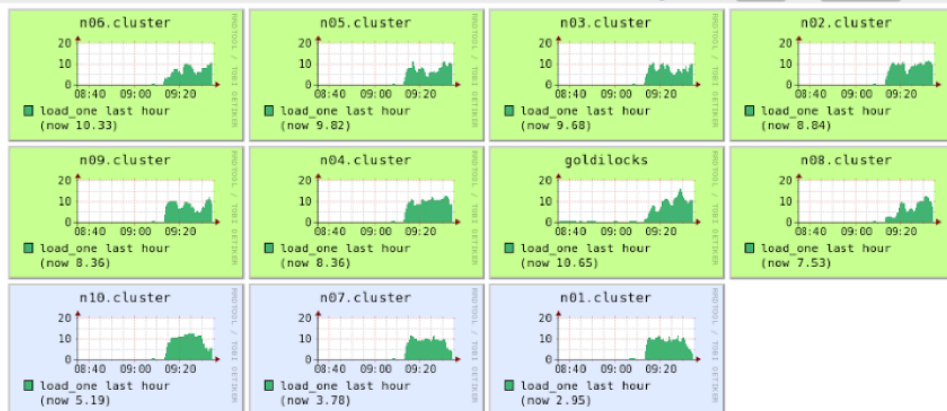
Avg Load (15, 5, 1m):
26%, 34%, 31%
Localtime:
2014-07-22 09:34



Overview of Goldenrod Cluster



Show Hosts: ☒ yes ☐ no | Goldenrod Cluster load_one last hour sorted descending | Columns Size



(Nodes colored by 1-minute load) | Legend

iostat

iostat -kxy 1 2

```
joshua@goldilocks:~> iostat -kxy 1 2
Linux 2.6.32-431.17.1.el6.x86_64 (goldilocks) 07/28/2014 _x86_64_ (32 CPU)

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.00    0.00   0.03   0.00    0.00   99.97

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s  avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.03    0.00   0.06   0.00    0.00   99.91

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s  avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     3.00     0.00    12.00     8.00     0.04    14.33   5.33   1.60
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     3.00     0.00    12.00     8.00     0.04    14.33   5.33   1.60
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

joshua@goldilocks ~> 
```


iostat

iostat -kxy 1 2

```
joshua@goldilocks:~> iostat -kxy 1 2
Linux 2.6.32-431.17.1.el6.x86_64 (goldilocks) 07/28/2014 _x86_64_ (32 CPU)

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.00    0.00   0.03   0.00    0.00   99.97

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s   avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.03    0.00   0.06   0.00    0.00   99.91

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s   avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     3.00    12.00    12.00     8.00     0.04    14.33   5.33   1.60
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     3.00    12.00    12.00     8.00     0.04    14.33   5.33   1.60
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

joshua@goldilocks ~>
```

kB read per second

iostat

iostat -kxy 1 2

```
joshua@goldilocks:~> iostat -kxy 1 2
Linux 2.6.32-431.17.1.el6.x86_64 (goldilocks) 07/28/2014 _x86_64_ (32 CPU)

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.00    0.00   0.03   0.00    0.00   99.97

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s  avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

avg-cpu:  %user   %nice %system %iowait  %steal   %idle
           0.03    0.00   0.06   0.00    0.00   99.91

Device:            rrqm/s   wrqm/s     r/s     w/s    rkB/s    kB/s  avgrq-sz  avgqu-sz   await  svctm   %util
sda                 0.00     0.00     0.00     3.00    12.00    12.00     8.00     0.04    14.33   5.33   1.60
sdb                 0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-0                0.00     0.00     0.00     3.00    12.00    12.00     8.00     0.04    14.33   5.33   1.60
dm-1                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00
dm-2                0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00     0.00   0.00   0.00

joshua@goldilocks ~>
```

kB read per second

kB written per second

Methods

Benchmarking

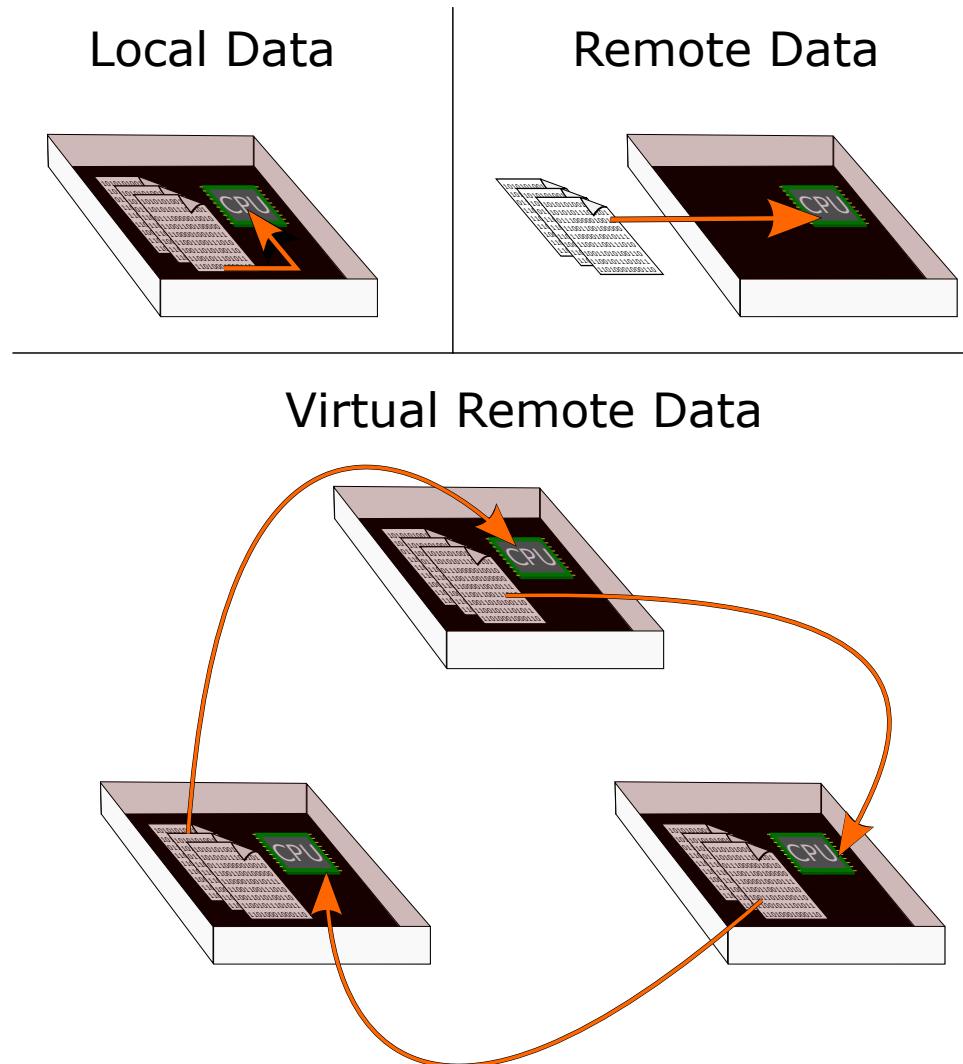
- bonnie++
- measure disk I/O

Hadoop jobs

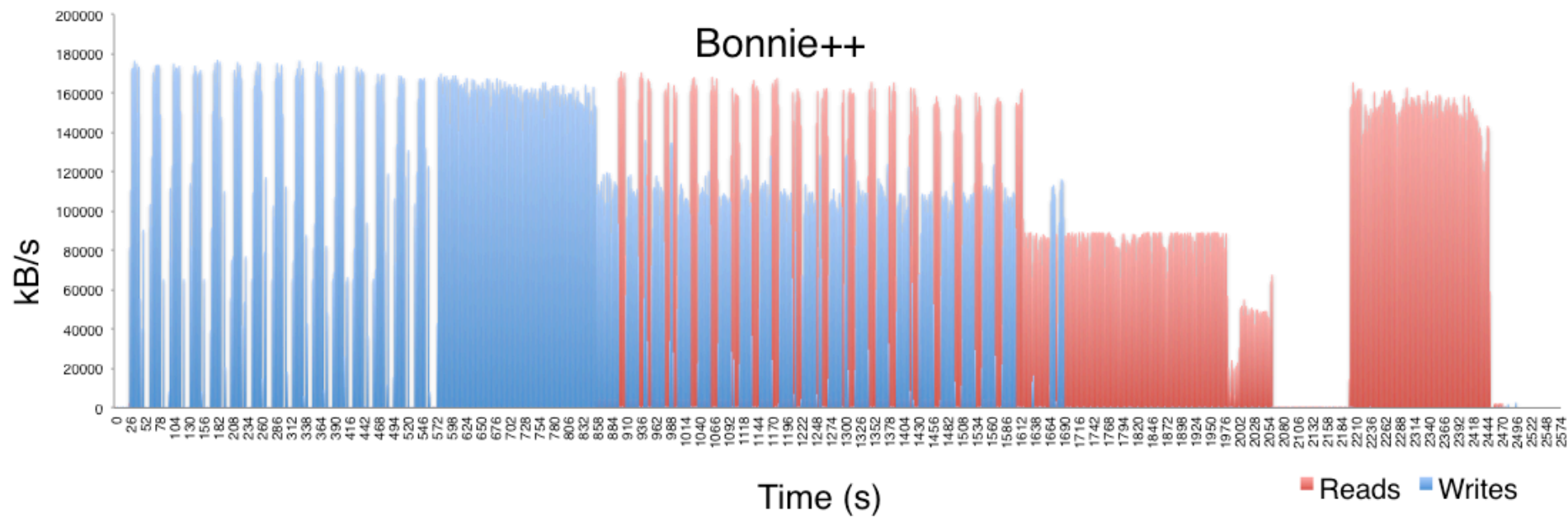
- teragen
- terasort

Hadoop jobs with remote data

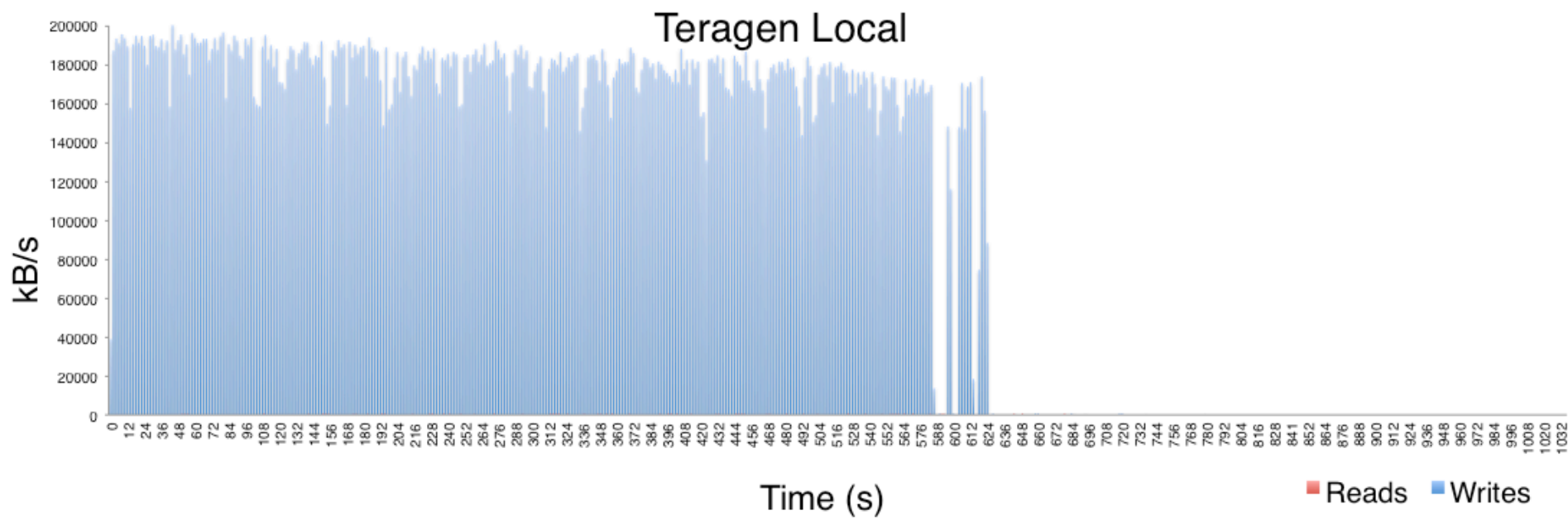
Methods



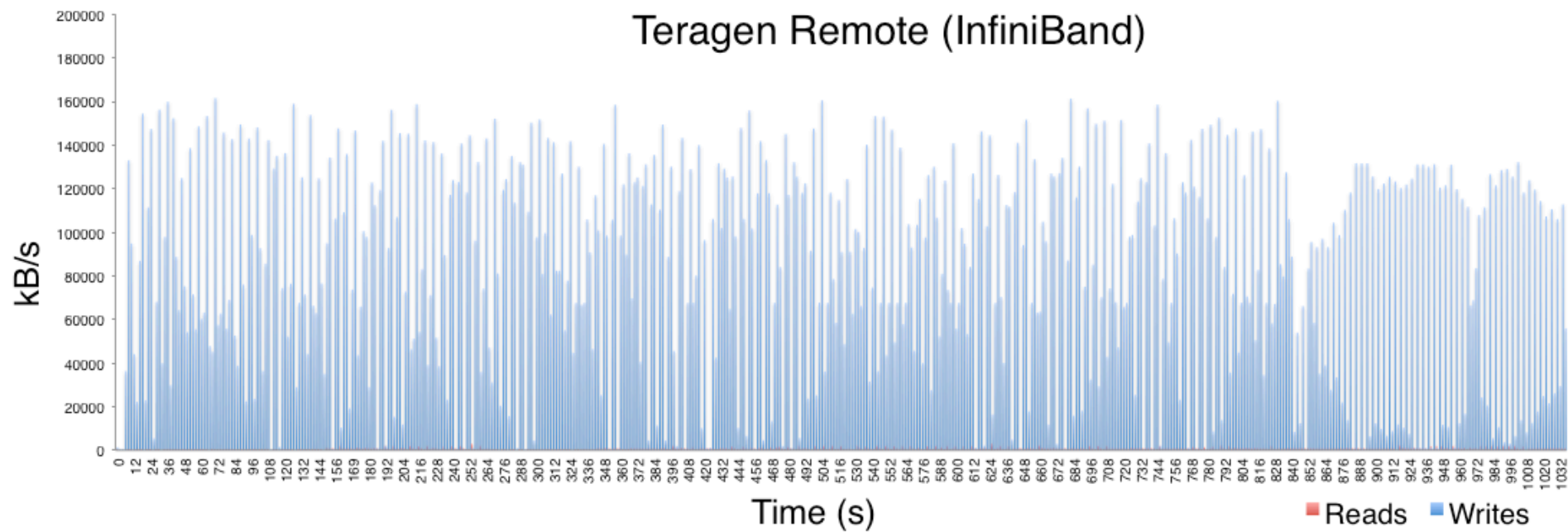
Results



Results

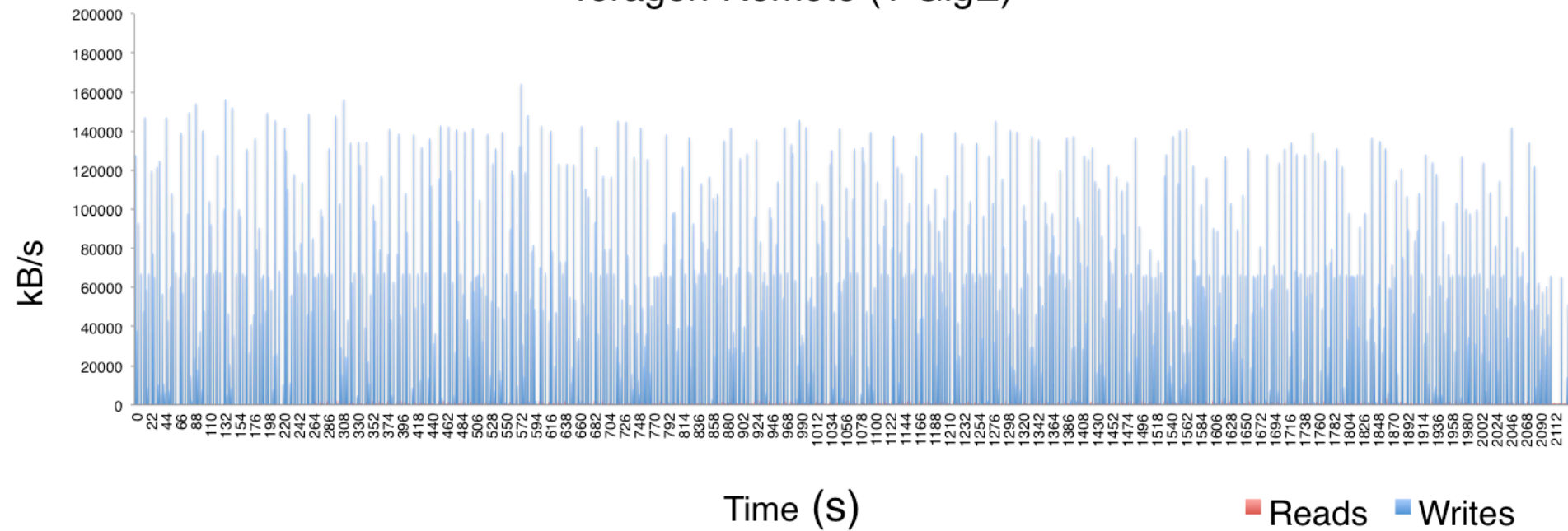


Results

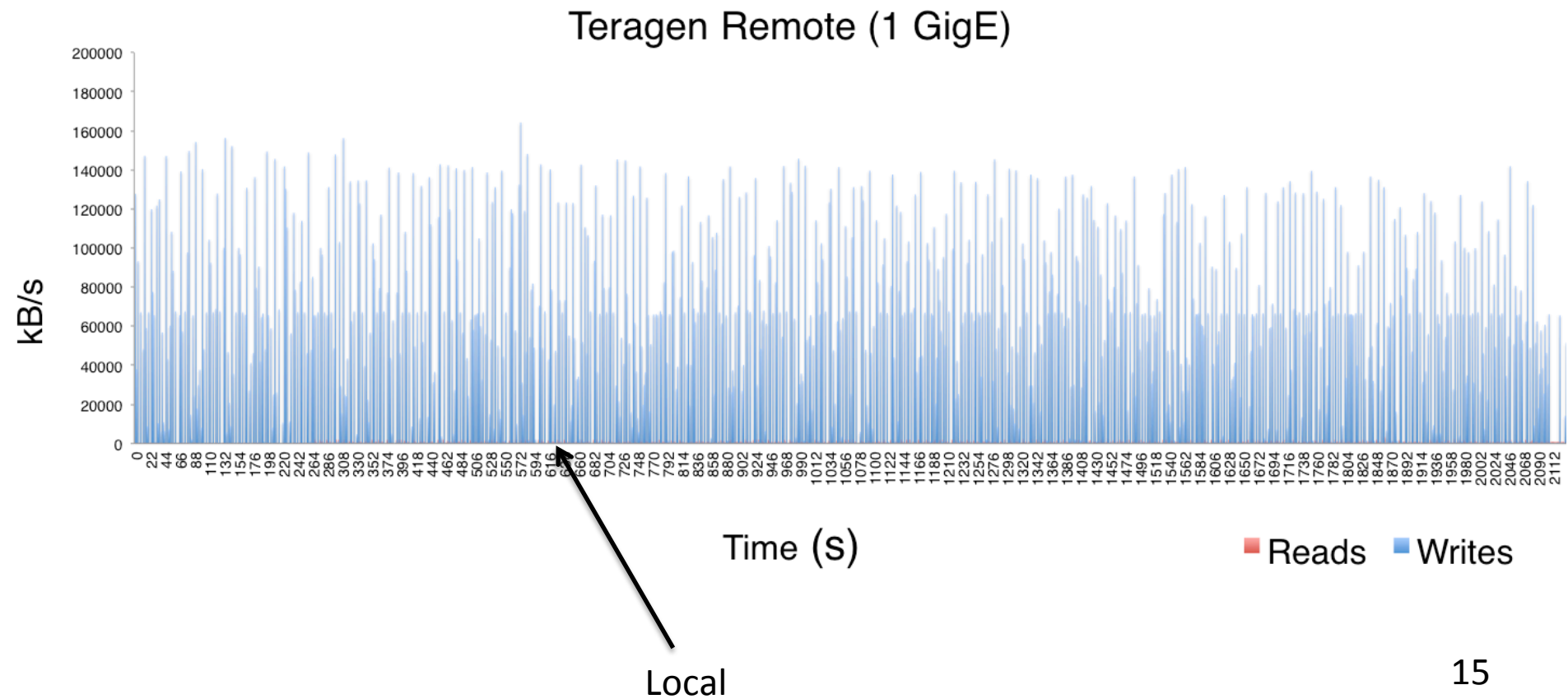


Results

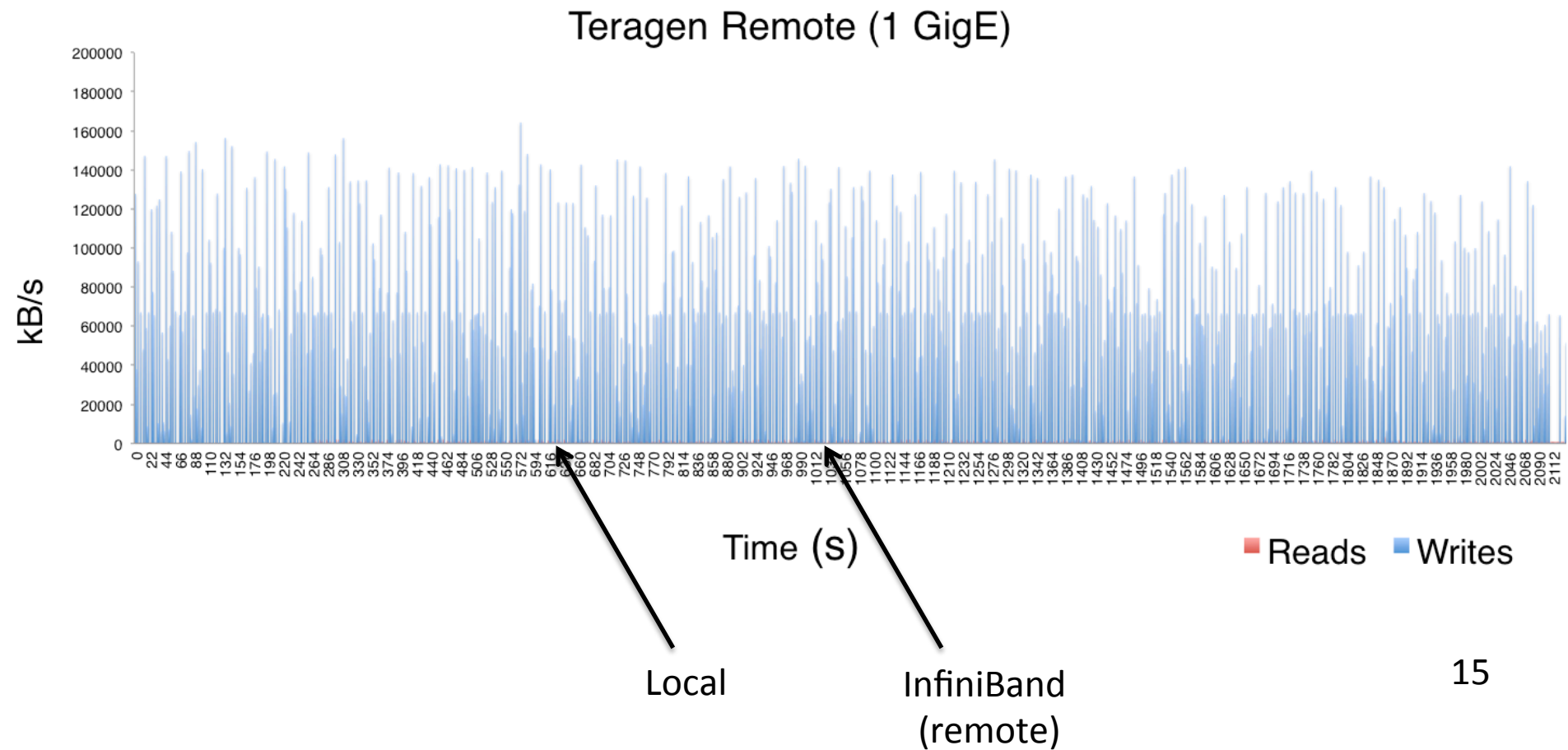
Teragen Remote (1 GigE)



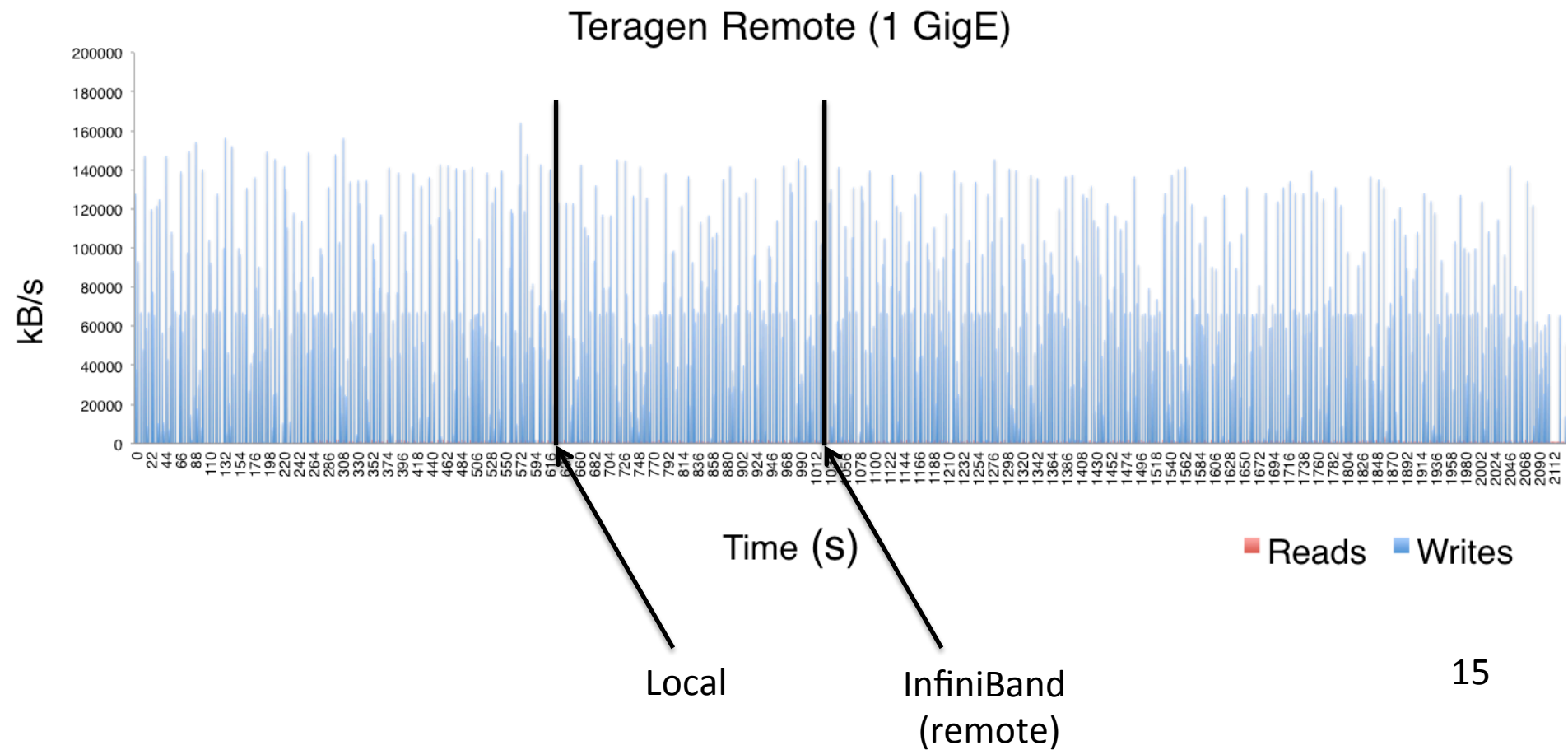
Results



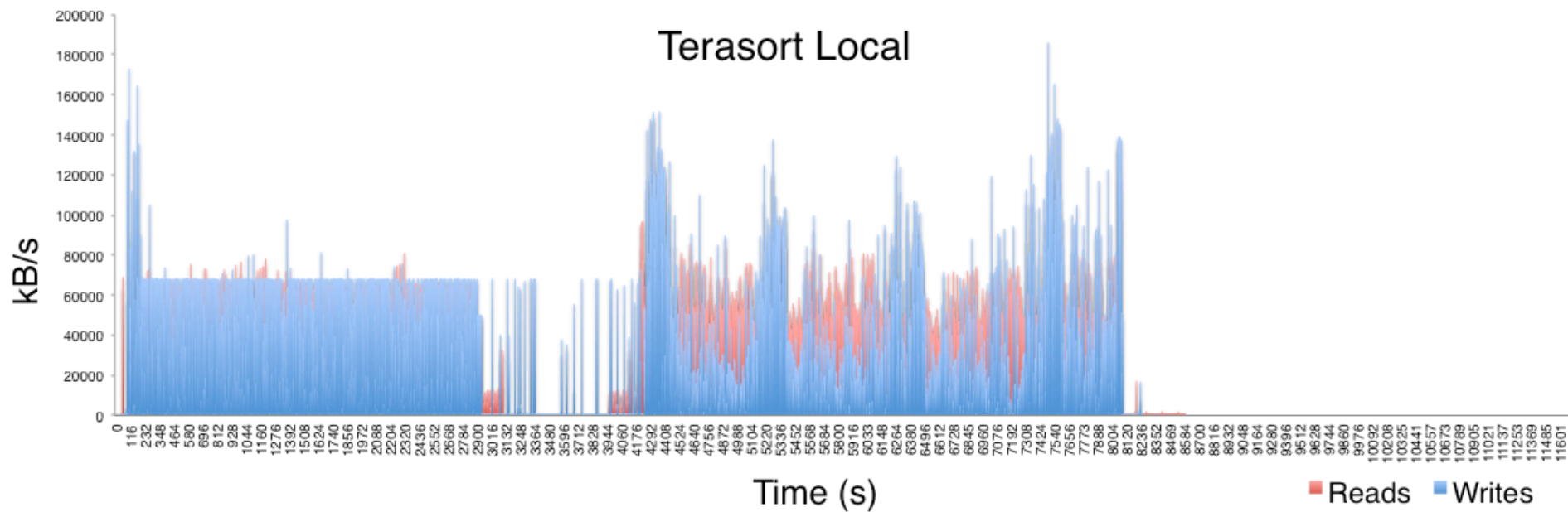
Results



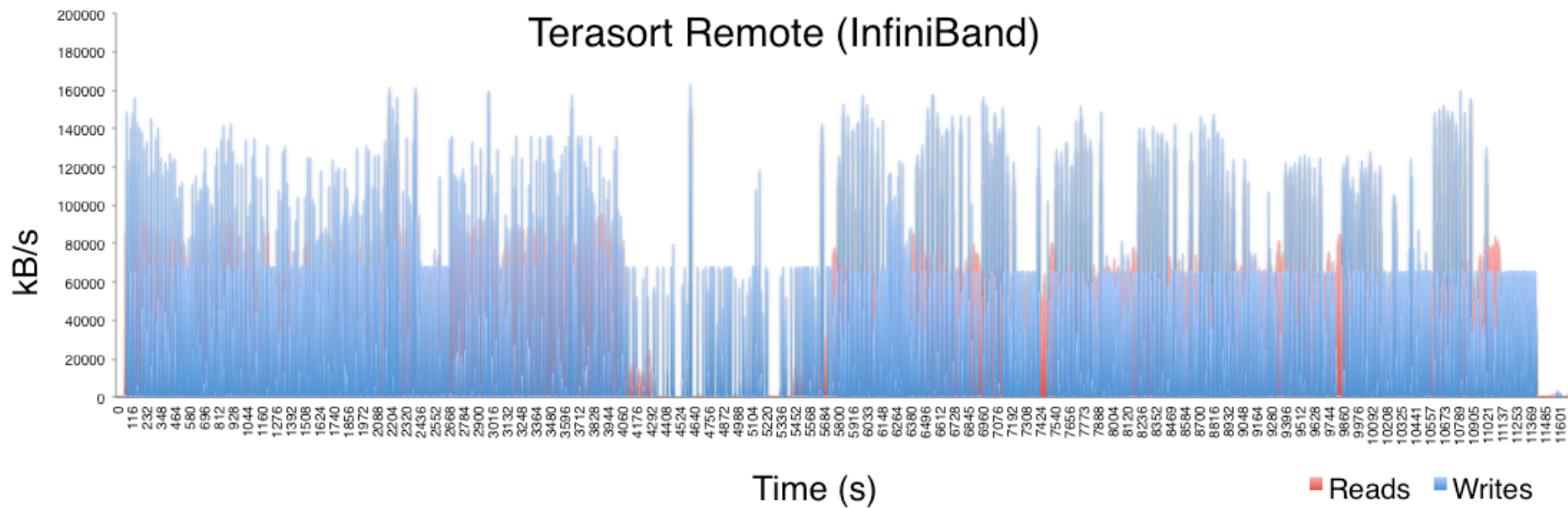
Results



Results



Results



Conclusion

Splunk

- Splunk app for HadoopOps is not suited to Hadoop MPv2/YARN

Ganglia

- Easy to configure and to extend

Effects of network latency

- Large impact when low connectivity
- Small, but noticeable impact for reasonable connectivity

Take-Aways and Successes

Monitoring I/O is easy (with the right tools)

- Successfully set up ganglia to monitor I/O
- Created visuals of I/O during Hadoop jobs

Benchmark of Hadoop jobs on local data and on remote data

- Performance suffers on data intensive jobs when data is stored remotely

Future Work

Write I/O monitoring application for Splunk

Evaluate effects of network latency with varying Hadoop parameters

- HDFS block size

Evaluating effects of network parameters

- Maximum transmission unit

Comparing performance on NFS to other file systems

Further examining trends in graphs

Questions?
/*Comments*/

